

Predicting Scientific Success Based on Coauthorship Networks

Emre Sarigöl, René Pfitzner*, Ingo
Scholtes, Antonios Garas, Frank Schweitzer

Chair of Systems Design, ETH Zurich

CH-8092 Zurich, Switzerland

*rpfitzner@ethz.ch

Abstract

We address the question to what extent the success of scientific articles is due to social influence. Analyzing a data set of over 100 000 publications from the field of Computer Science, we study how centrality in the coauthorship network differs between authors who have highly cited papers and those who do not. We further show that a machine learning classifier, based only on coauthorship network centrality measures at time of publication, is able to predict with high precision whether an article will be highly cited five years after publication. By this we provide quantitative insight into the social dimension of scientific publishing – challenging the perception of citations as an objective, socially unbiased measure of scientific success.

1 Introduction

Quantitative measures are increasingly used to evaluate the performance of research institutions, departments, and individual scientists. Measures like the absolute or relative number of published research articles are frequently applied to quantify the *productivity* of scientists. To measure the *impact* of research, citation-based measures like the total number of citations, the number of citations per published article or the h-index [9], have been proposed. Proponents of such citation-based measures or rankings argue that they allow to quantitatively and objectively assess the *quality* of research, thus encouraging their use as simple proxies for the *success* of scientists, institutions or even whole research fields. The intriguing idea that by means of citation metrics the task of assessing research quality can be “outsourced” to the *collective intelligence* of the scientific community, has resulted in citation-based measures becoming increasingly popular among research administrations and governmental decision makers. As a result, such measures are used as one criterion in the evaluation of grant proposals and research institutes or in hiring committees for faculty positions. Considering the potential impact for the careers of - especially young - scientists, it is reasonable to take a step back and ask a simple question: To what extent do *social factors* influence the number of citations of their articles? Arguably, this question challenges the perception of science as a systematic pursuit for objective truth, which ideally should

be free of personal beliefs, biases or social influence. On the other hand, quoting Werner Heisenberg [8], “*science is done by humans*”, it would be surprising if specifically scientific activities were free from the influences of social aspects. Whereas often the term “social influence” has a negative connotation, we don’t think that social influence in science necessarily stems from malicious or unethical behavior, like e.g. nepotism, prejudicial judgments, discrimination or in-group favoritism. We rather suspect that, as a response to the increasing amount of published research articles and our limited ability to keep track of potentially relevant works, a growing importance of social factors in citation behavior is due to natural mechanisms of *social cognition* and *social information filtering*.

In this paper we address this issue by studying the influence of social structures on scholarly citation behavior. Using a data set comprising more than 100 000 scholarly publications by more than 160 000 authors, we extract time-evolving coauthorship networks and utilize them as a proxy for the evolving social network of the scientific discipline *computer science*. Based on the assumption that the centrality of scientists in the resulting social network is indicative for the *visibility* of their work, we then study to what extent the “*success*” of research articles in terms of citations can be predicted using only knowledge about the embedding of authors in the social network *at time of publication*. Our prediction method is based on a random forest classifier and utilizes a set of complementary network centrality measures. We find strong evidence for our hypothesis that authors whose papers are highly cited in the future have - on average - a significantly higher centrality in the social network at the time of publication. Remarkably, we are able to predict whether an article will belong to the 10% most cited articles with a precision of 60%. We argue that this result quantifies the existence of a *social bias*, manifesting itself in terms of visibility and attention, and influencing measurable citation “success” of researchers. The presence of such a social bias not only highlights problems with current publication and citation practices. It also threatens the interpretation of citations as *objectively awarded esteem*, which is the justification for using citation-based measures as universal proxies of *quality* and *success*.

The remainder of this article is structured as follows: In section 2 we review a number of works that have studied scientific collaboration structures as well as their relation to citation behavior. In section 3 we describe our data set and provide details of how we construct time-evolving coauthorship networks. We further introduce a set of network-theoretical measures which we utilize to quantitatively assess the centrality and embedding of authors in the evolving coauthorship network. In section 4 we introduce a number of hypotheses about the relations between the position of authors in the coauthorship network and the future success of their publications. We test these hypotheses and obtain a set of candidate measures which are the basis for our prediction method described in section 5. We summarize and interpret our findings in section 6 and discuss their implications for the application of citation-based measures in the quantitative assessment

of research.

2 The Complex Character of Citations

It is remarkable that, even though citation-based measures have been used to quantify research impact since almost sixty years [6], a complete *theory of citations* is still missing. In particular, researchers studying the social processes of science have long been arguing that citations have different, complex functions that go well beyond a mere attribution of credit [14]. At the level of scientific articles, a citation can be interpreted as a “discursive relation”, while at the level of authors citations have an additional meaning as expression of “professional relations” [14]. Additional interpretations have been identified at aggregate levels, like e.g. social groups, institutions, scientific communities or even countries citing each other. These findings suggest that citations are indeed a complex phenomenon which have both cognitive and a social dimension [14, 20]. This questions an oversimplified interpretation of citations as objective quality indicator. The complex character of scholarly citations was further emphasized recently [13]. Here, the authors argue that, apart from an attribution of scientific merit, references in scientific literature often serve as a tool to guide and orient the reader, to simplify scientific writing and to associate the work with a particular scientific community. Furthermore, they highlight that citation numbers of articles are crucially influenced not only by the popularity of a research topic and the size of the scientific community, but also by the number of authors as well as their prominence and visibility.

Facilitated by the wide-spread availability of scholarly citation databases, some advances in the understanding of the dynamics of citations have been made in the last years. Generally, citation practices seem to differ significantly across different scientific disciplines, and thus complicating the definition of universal citation-based impact measures. However, the remarkable finding that – independent of discipline – citations follow a log-normal distribution and can be rescaled in such a way that citation numbers become comparable [21, 22], suggests that the mechanisms behind citation practices are universal across disciplines, and differences are mainly due to differing community sizes.

Additionally to investigations of the differences across scientific communities, the relations between citations and coauthorships were studied in recent works. Using data from a number of scientific journals, it was shown that the citation count of an article is correlated both with the number of authors and the number of institutions involved in its production [5, 12]. Studying data from eight highly ranked scientific journals, it was shown [11] that a) single author publications consistently received the lowest number of citations and b) publications with less than five coauthors received less citations than the average article. Studying citations between individuals rather than articles, in [16] it was observed that coauthors tend to cite each other sooner after the

publication of a paper (compared to non-coauthors). Further, the authors showed that a strong tendency towards reciprocal citation patterns exists. Although these findings already indicate that social aspects influence citing behavior, in this work we are going to quantitatively reveal the true extent of this influence.

Going beyond a mere study of direct coauthorship relations, first attempts to study *both* citation and coauthorship structures from a *network perspective* have been made recently. Aiming at a measure that captures both the *amount* as well as the *reach* of citations in a scientific community, a citation index that incorporates the distance of citing authors in the collaboration network was proposed [2]. Another recent study [23] used the topological distance between citing authors in the coauthorship network to extend the notion of self-citations. Interestingly, apart from direct self-citations, this study could not find a strong tendency to cite authors that are close in the coauthorship network.

Different from previous works, in this article we study correlations between the *centrality* of authors in collaboration networks and the *citation success* of their research articles. By this we particularly extend previous works that use a network perspective on coauthorship structures and citation patterns. Stressing the fact that *social relations* of authors play an important role for how much attention and recognition their research receives, we further contribute a quantitative view on previously hypothesized relations between the *visibility* of authors and citation patterns.

3 Time-Evolving Collaboration and Citation Networks

In this work we analyze a data set of scholarly citations and collaborations obtained from the Microsoft Academic Search ¹ (MSAS) service. The MSAS is a scholarly database containing more than 35 Million publication records from 15 scientific disciplines. Using the Application Programming Interface (API) of this service, we extracted a subset of more than 100 000 computer science articles, published between 1996 and 2008, in the following way: First, we retrieved unique numerical identifiers (IDs) of the 20 000 highest ranked authors in the field of *computer science*. This ranking is the result of an MSAS internal “field rating”, taking into account several scholarly metrics of an author (number of publications, citations, h-index) and comparing them to the typical values of these metrics within a certain research field. As the goal was to build coauthorship and citation networks of reasonable size, in a second step we chose 1000 authors i.i.d. uniformly from the set of these 20 000 authors. In the third step, we obtained information on coauthors, publication date, as well as the list and publication date of citing works for all the publications authored by these 1000 authors between 1996 and 2008. This results in a data set consisting of a total of 108 758 publications from the field computer science, coauthored by a total

¹<http://academic.research.microsoft.com>

of 160 891 researchers. Each publication record contains a list of author IDs, which, by means of disambiguation heuristics internally applied by the MSAS service, uniquely identify authors independent of name spelling variations. The absence of name ambiguities is one feature that sets this data set apart from other data sets on scholarly publications that are used frequently. Based on this data set we extracted a *coauthorship network*, where nodes represent authors and links represent coauthorship relations between authors. In addition, using the information about citing papers, we extracted *citation dynamics*, i.e. the time evolution of the number of citations of all publications in our data set. Similar to earlier works, we argue that the coauthorship network can be considered a first-order approximation of the complete scientific collaboration network [16]. Based on the publication date of an article, we additionally assign time stamps to the extracted coauthor links – thus obtaining time-evolving coauthorship networks.

We analyze the evolution of the coauthorship network using a sliding window of two years in which we aggregate all coauthorships occurring within that time. Starting with 1996, we slide this window in one year increments and obtain a total of 11 time slices representing the evolution of collaboration structures between 1996 and 2008. We use an extended time-window of two years to account for the continuing effect of a coauthorship in terms of awareness about the coauthors works. Although larger time windows are certainly possible (and their effects interesting to investigate), in this work we are less concerned with the optimal time-window size and consistently use the above described approach. However, performed consistency checks with varying time-window sizes suggest robustness of our results.

Table 1 summarizes the number of nodes and links in the coauthorship network, the number of publications in each time slice as well as the fractional size of the largest connected component (LCC). Note that the time-aggregated network (overall) forms one giant component with only a minor fraction of isolated nodes, whereas some of the time slices fall apart into many separated components. Note also that the size of the largest connected component is increasing with time, which may indicate either a possible bias in the coverage of the MSAS database to favour newer articles, or an increase of “collaborativeness” in science. As we are going to perform a social network analysis of the collaboration time slices – and some measures (like eigenvector centrality) are not well-defined for unconnected graphs – we apply all of the following analysis always on the largest connected component. For each network corresponding to one two-year time slice, we compute a number of node-level metrics that allow us to quantitatively monitor the evolution of network positions for all authors. In particular, we compute *degree centrality*, *eigenvector centrality*, *betweenness centrality* and *k-core centrality* of authors. For details on the used centrality measures, please refer to the Supplementary Material or the textbook by Newman [18]. Here we utilize implementations of these measures provided by the igraph package [4].

A major focus of our work is to assess the predictive power of an author’s position in the coauthorship network for the citation success of her future articles. To do so we adopt a so

Year	LCC fraction	Links	Nodes	Publications
1996-1997	0.18	61 046	2845	1160
1997-1998	0.37	130 938	6381	3070
1998-1999	0.45	153 412	8470	4054
1999-2000	0.50	186 318	10 413	5320
2000-2001	0.60	358 188	13 451	6561
2001-2002	0.63	413 846	15 309	7026
2002-2003	0.74	542 912	20 238	9193
2003-2004	0.77	653 224	23 624	10 608
2004-2005	0.79	745 352	26 258	11 430
2005-2006	0.83	889 996	29 886	12 919
2006-2007	0.84	914 614	32 412	13 568
2007-2008	0.86	858 554	35 255	14 214
Overall	0.99	5 324 330	160 891	108 758

Table 1: Number of papers and size of the collaboration network 2-year subgraphs between 1995-2008 used in our study.

called *hindcasting approach*: For each publication p published in a given year t , we extract the list of coauthors as well as the LCC of the coauthorship network in the time slice $[t - 2, t]$, and calculate the centrality measures. Based on the citation data, we furthermore calculate the number of citations c_p paper p gained within a time frame of *five years* after publication, i.e. in the time slice $[t, t + 5]$.

In particular, we are interested in those publications that are among the most successful ones. Defining *success* is generally an ambiguous endeavor. As justified in the introduction, here we take the (controversial) viewpoint that success is directly measurable in number of citations. We specifically focus on a very simple notion of success in terms of *highly cited papers* and, similar to [17], assume that a paper is *successful* if five years after publication it has more citations than 90% of all papers published in the same year. We refer to the set of successful papers in year t as $P_{\uparrow}(t)$. The set of remaining papers, i.e. those published at time t that are cited less frequently than the top 10%, is denoted as $P_{\downarrow}(t)$.

4 Statistical Dependence of Coauthorship Structures and Citations

Having a large social network and “knowing the right people” often is a prerequisite for career success. However, science is often thought to be one of the few fields of human endeavor where success depends on the quality of an authors’ work, rather than on her social connectedness. Given the time evolving coauthorship network, as well as the observed success (or lack thereof) of a publication, we investigate two research questions, aiming to quantify the aspect of social influence on citation success. First, we examine whether there is a general statistical dependency of central authors in the coauthorship network to publish papers that are more successful than non-central. Second, we investigate whether the inverse effect is present and the success of a paper influences the future coauthorship centrality of its authors.

4.1 Effects of Author Centrality on Citation Success

To quantify the first research question we test the following hypothesis.

H1: *At the time of publication, authors of papers in $P_{\uparrow}(t)$ are more central in the coauthorship network than authors of articles in $P_{\downarrow}(t)$.*

As papers often have more than one author, for each paper we will consider only the coauthorship network centralities of the author with the highest coauthorship degree, and refer to this as the *coauthorship centrality of the paper*. This choice is motivated by the intuition that the centrality of the best connected coauthor should provide the major amount of (socially triggered) visibility for the publication. We test **H1** by comparing coauthorship centrality distributions of papers in $P_{\uparrow}(t)$ and $P_{\downarrow}(t)$ for each year t . In order to compare the centrality distributions, we apply a *Wilcoxon-Mann-Whitney two-sample rank-sum-test* [15]. For each of the four centrality metrics we test the null hypothesis that coauthorship centrality distributions of papers in $P_{\uparrow}(t)$ and $P_{\downarrow}(t)$ are the same against the alternative hypothesis that the centrality distribution of papers in $P_{\uparrow}(t)$ is stochastically larger than the one of papers in $P_{\downarrow}(t)$. The p-values of the tests as well as the corresponding averages and variances of the four considered centrality metrics in the two sets are shown in Table 2. For all considered centrality metrics p-values are well below a significance level of 0.01. This leads us to safely reject the null hypothesis, concluding that coauthorship centrality metrics of papers in $P_{\uparrow}(t)$ are stochastically larger than of those papers in $P_{\downarrow}(t)$. This result indicates that all considered centrality metrics in the coauthorship network, at the time of publication of a paper, are indicative for future paper success. Note however, that this statistical dependency is more complicated than the linear Pearson or the more general Spearman correlation. Indeed, all the considered social network metrics are only weakly, if at all, correlated with citation numbers (see Supplementary Material). To what

	p-value	$\langle P_{\downarrow} \rangle$	$\langle P_{\uparrow} \rangle$	$\text{var } P_{\downarrow}$	$\text{var } P_{\uparrow}$
<i>k</i>-core	1.28×10^{-115}	3.33×10^1	4.39×10^1	1.20×10^4	7.18×10^3
Eigenvector	2.52×10^{-34}	2.87×10^{-3}	5.60×10^{-4}	2.58×10^{-3}	5.40×10^{-4}
Betweenness	1.19×10^{-68}	6.30×10^5	1.52×10^6	4.19×10^{12}	1.58×10^{13}
Degree	5.63×10^{-125}	9.38×10^1	1.57×10^2	1.02×10^5	1.13×10^5

Table 2: P-values of one sided Wilcoxon-Mann-Whitney test. This quantifies whether the distribution of centralities of authors of articles in P_{\uparrow} are (in a statistical sense) larger than those of authors of articles in P_{\downarrow} . Also shown are the means and variances of the centrality metrics in the two sets.

	Top 10%	Top 5%	Top 2%	Top 1%
<i>k</i>-core	0.22 0.21	0.17 0.16	0.07 0.07	0.01 0.01
Eigenvector	0.11 0.11	0.06 0.06	0.01 0.01	0.01 0.01
Betweenness	0.20 0.20	0.13 0.13	0.11 0.11	0.11 0.11
Degree	0.20 0.20	0.15 0.15	0.10 0.09	0.07 0.07
Intersection	0.36 0.15	0.27 0.11	0.17 0.06	0.12 0.04
# papers	3700	1844	730	362

Table 3: First table entry indicates what fraction of papers, that have authors which are within the set of author with Top $x\%$ centrality metrics, are also Top $x\%$ of all papers in terms of citation success ($P(\text{toppaper}|\text{topmetric})$). Second table entry indicates what fraction of papers, that are Top $x\%$ of all papers in terms of citation success, have authors which are within the set of author with Top $x\%$ centrality metrics ($P(\text{topmetric}|\text{toppaper})$). Row *Intersection* indicates the intersection of all the above considered centrality metrics.

extent citation success and coauthorship network centrality are statistically dependent is summarized in Table 3. Left entry of each cell indicates what fraction of papers, that have authors with Top $x\%$ centrality metrics, belong to the Top $x\%$ of all papers in terms of citation success ($P(\text{toppaper}|\text{topmetric})$). Right entry of each cell indicates what fraction of papers, that are Top $x\%$ of all papers in terms of citation success, have authors which are within the set of authors with Top $x\%$ centrality metrics ($P(\text{topmetric}|\text{toppaper})$). From these results, we conclude two observations: First, the probabilities in every cell are well below 1, indicating the absence of a simple linear (Pearson) correlation. Second, especially considering k -core centrality, knowing a paper is Top 10% successful, the conditional probability that it was written by an author with Top 10% k -core centrality, is $P(\text{topmetric}|\text{toppaper}) = 0.21$. Additionally, Table 3 indicates that vice versa $P(\text{toppaper}|\text{topmetric}) = 0.22$ of all papers that are published by authors with Top 10% k -core centrality, are successful. Considering the intersection of all four centrality metrics, we even find that $P(\text{toppaper}|\text{topmetric}) = 0.36$ of all papers published by the Top 10% central

authors in all four coauthorship centrality metrics, are within the Top 10% cited papers. We will use this observation as basis for a naive Bayes classifier in section 5.

4.2 Coevolution of Coauthorship and Citation Success

In the previous section we studied the question whether the centrality of authors in the coauthorship network is indicative for the success of publications in terms of citations. Our results suggested that centrality in coauthorship networks is indeed indicative for citation success. In the following we study the inverse relation and ask whether a shift in citation success of an author is indicative for her future position in the coauthorship network. To answer this question, we consider all authors who published an article both at time t and five years later at $t + 5$. We then categorize them based on the citation success of their articles published at time t and time $t + 5$. We introduce two sets of authors: Set $A_{\searrow}(t)$, which is the set of those authors who at time t had at least one publication in class $P_{\uparrow}(t)$, but who at time $t + 5$ did not have an article in class $P_{\uparrow}(t + 5)$ anymore. Set $A_{\nearrow}(t)$ contains all authors who at time t had no article in class $P_{\uparrow}(t)$ but who at time $t + 5$ published at least one article that falls in class $P_{\uparrow}(t + 5)$. In addition, we again record the coauthorship centralities of authors in these two sets for time windows $[t - 2, 2]$ and $[t + 3, t + 5]$.

For authors of set A_{\nearrow} we test the following hypothesis:

H2: *Authors that experience a positive shift in their citation success (i.e. authors in A_{\nearrow}) will become more central in the coauthorship network.*

Complementary to **H2**, for authors in set A_{\searrow} we hypothesize:

H3: *Authors that experience a negative shift in their citation success (i.e. authors in A_{\searrow}) will become less central in the coauthorship network.*

In order to test for **H2** and **H3**, we apply a *pairwise Wilcoxon-Mann-Whitney* test. To verify **H2** we test if the centralities of authors have decreased in the case of a decrease in publication success from time t to $t + 5$. To verify **H3** we test if the centralities of authors have increased in the case of an increase in publication success from time t to $t + 5$. Results of these hypotheses tests are presented in Table 4. For authors in A_{\nearrow} we observe p -values much lower than the 0.01 significance threshold of the alternative hypothesis testing (**H2**). We hence find evidence that authors in A_{\nearrow} experience a significant increase in k -core, betweenness and degree centrality. Reversely, results for authors in A_{\searrow} suggest a significant drop in k -core, eigenvector and degree centrality. Based on these results we cannot reject hypothesis **H2** and **H3**, indicating that there is significant influence of author citation success on her future coauthorship network centrality.

As an illustration of citation and coauthorship dynamics, Figure 1 shows part of the coauthorship network. Color intensity of the nodes is scaled to their degree centrality, while node size is scaled

centrality measure & alternative	A_{\searrow}	A_{\nearrow}
$k\text{-core}(t) > k\text{-core}(t+5)$	3.15×10^{-11}	1
$k\text{-core}(t) < k\text{-core}(t+5)$	1	3.04×10^{-55}
$\text{ev-centr}(t) > \text{ev-centr}(t+5)$	5.18×10^{-14}	0.86
$\text{ev-centr}(t) < \text{ev-centr}(t+5)$	1	0.14
$\text{bw-centr}(t) > \text{bw-centr}(t+5)$	0.23	1
$\text{bw-centr}(t) < \text{bw-centr}(t+5)$	0.77	7.29×10^{-30}
$\text{degree}(t) > \text{degree}(t+5)$	6.69×10^{-11}	1
$\text{degree}(t) < \text{degree}(t+5)$	1	7.72×10^{-62}
# authors	521	648

Table 4: P-values of Wilcoxon-Mann-Whitney test for different coauthorship centralities and alternative hypotheses. Column A_{\searrow} presents p-values for authors in set A_{\searrow} , column A_{\nearrow} presents p-values for authors in set A_{\nearrow} .

to their betweenness centrality. A very strong community structure is clearly visible. Furthermore, we highlighted in red one particular author that belonged to group $A_{\nearrow}(t)$, i.e. authors who did not have a paper in P_{\uparrow} in 2002, but did so in 2007. Thus, in the considered five year span the highlighted author moved from a position in the periphery of the coauthorship network to a position in the center. Not only the authors' degree centrality increased (see size of the node as well as joined red-colored links), but also betweenness centrality improved highly.

Note that already in 2002 the author had comparatively high betweenness and degree centrality, which –according to our previous discussion– provided an ideal starting point for citation success in 2007.

5 Predicting Successful Publications

In the previous sections we presented evidence for the existence of statistical dependencies between authors' coauthorship centrality and the success of their publications. Results suggested that several coauthorship centrality metrics are indicative for citation success. However, we did not identify one single such centrality metric, especially we did not find that the mere number of coauthors is sufficient for a paper to become highly cited. Instead, this seems to be dependent on more than one network measure. In this section we present a machine learning classifier to predict whether a publication will be highly cited, based on several features of the authors position in the coauthorship network.

Previous works have already attempted to predict citation success. For example in [10], the predictive power of the past h-index for the future h-index of a scientist was presented. Furthermore, in [1] additional indicators like, e.g. the length of the career or the number of articles in certain journals, have been integrated into a model to predict the future h-index of scientists. The au-

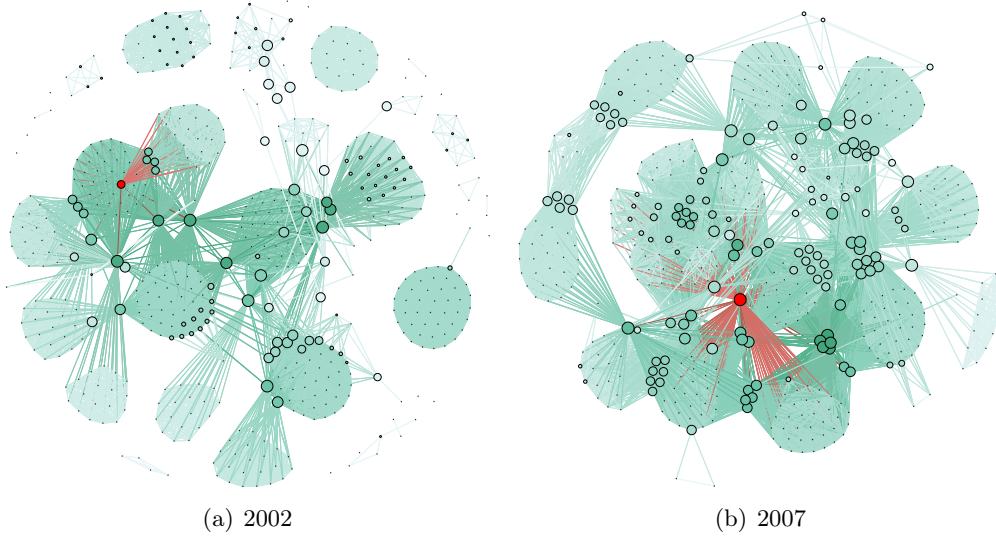


Figure 1: Illustration of correlation between citation success and centrality in the coauthorship network. Color intensity of the nodes is scaled according to their degree centrality and size of nodes is scaled according to their betweenness centrality.

thors of [17] compare the number of citations an article has received at a given point in time with the expected value in a preferential attachment model for the citation network. Deriving a z-score, the authors present a prediction of which papers will be highly cited in the future. Recently the authors reevaluate their earlier predictions and confirm the predictive power of their approach [19]. Whereas these three approaches attempt to predict success based on past citation dynamics, they do not investigate the underlying mechanisms that lead to citation success. Here we address this fundamental question and try to predict citation success based merely on coauthorship network centrality of authors. Clearly, many different factors will contribute to scientific success. In this work, however, we focus on the social component (based on the coauthorship network) in order to highlight the influence of social, not necessarily merit-based, mechanisms on publication success.

Based on the observed relations between author centralities in the coauthorship network and the success of their publications presented in section 4, in this section we investigate whether we can predict a paper’s future success. In particular, we try to predict whether a paper will be highly cited five years after its publication based on measures of author centrality in the coauthorship network.

In section 4.1 we presented insights about the statistical dependency of citation success and several social network centrality measures (see Table 3). These results suggest that a naive Bayes predictor for citation success can already yield quite useful results, predicting whether or not a

paper will be **toppaper**, given ex ante knowledge about **topmetric** of the authors. Using k -core centrality as a basis, we apply the following classification rule:

If a paper is authored by a top 10% k -core centrality author, then the paper will be among the top 10% most cited papers five years after publication.

To evaluate the goodness of this prediction, we will consider the error measures *precision* and *recall*². Observing that for k -core centrality in a 10% success scenario it is $P(\text{topmetric}|\text{toppaper}) = 0.21\%$ as well as $P(\text{toppaper}|\text{topmetric}) = 0.22\%$ and the fact that for a naive Bayes classifier $\text{recall} = P(\text{topmetric}|\text{toppaper})$ and $\text{precision} = P(\text{toppaper}|\text{topmetric})$ holds, one sees that a classifier with the above rule yields $\text{recall} = 21\%$ and $\text{precision} = 22\%$. Similarly, instead of k -core centrality other network measures presented in Table 3 can be used as basis for the above classification rule. As earlier works have tried to predict the success of papers based on the number of coauthors [11], using degree centrality as basis for the above classification rule directly extends these attempts, yielding $\text{recall} = 20\%$ and $\text{precision} = 20\%$. Note, however, that degree centrality accumulates all coauthorships that have been established within the two-year sliding window of our analysis, not just the coauthorships of the paper under consideration.

We now ask whether a multi-dimensional naive Bayes classifier can improve this single metric classification result. Taking into account the intersection of all considered centrality metrics, we consider the following classification rule:

If a paper is authored by an author with a top 10% betweenness centrality, degree centrality, k -core centrality and eigenvector centrality, then the paper will be among the top 10% most cited papers five years after publication.

Using this classifier, we achieve even better classification of $\text{precision} = 0.36\%$, however diminishing recall to $\text{recall} = 0.15\%$. Whereas these results already show that a naive Bayes classifier can yield interesting insights, in the following we will present a more sophisticated Machine Learning approach, taking multiple network centrality features into account and improving classification errors.

We first construct a feature vector for every publication as follows. For each publication appearing in year t , we extract all coauthors and compute the maximum and minimum of their centralities in the coauthorship network constructed based on the time window $[t-2, t]$. Then, for each publication we build a feature vector with 10 features containing the maximum and minimum of the centrality metrics considered earlier (*degree*, *eigenvector*, *betweenness* and *k -core*), as well as the number of coauthors and the cumulative number of authors a paper has referenced. We then classify all publications regarding whether they fall in P_{\uparrow} or P_{\downarrow} according to the aforementioned

²See Supplementary Material for a general definition of precision and recall

Nr.Publications	Precision	Recall	F-Score
36000	60%	18%	28%

Table 5: Error estimates of the Random Forest classifier to predict success of papers.

publication classes, with P_{\uparrow} defined as the set of the top 10% cited publications and P_{\downarrow} as the remaining 90%.

The classification is done using a Random Forest classifier [3], extending the concept of classification trees³. In general, the Random Forest is known to yield accurate classifications for data with a large number of features [3]. Furthermore, it is a highly scalable classification algorithm, eliminating the need for separate cross validation and error estimation, as these procedures are part of the internal classification routine.⁴

Table 5 summarizes precision, recall, and F-score of the resulting classification. Comparing this result with the expectation from a random guess, which will correctly pick one of the top 10% publications only in 10% of the cases, the achieved precision of 60% is striking. In particular, by only considering positional features of authors in the coauthorship network, we are able to achieve *an increase of factor six in predictive power* compared to a random guess. Also, we obtain a *recall* value of 18%, meaning that our classifier correctly identified about one fifth of all of the top 10% papers in a given research field. As a random guess would yield a recall of 10%, the Random Forest classifier *improves recall by 80%*.

This result allows for two conclusions: First, the fact that a high-dimensional random forest classifier performs better than a naive Bayes classifier, makes clear that social influence on scientific success cannot be measured by a single value. Second, and most importantly, that by *solely considering metrics of social influence*, such a classifier is able to predict scientific success with high precision.

6 Discussion and Conclusions

Using a data set on more than 100 000 scholarly publications authored by more than 160 000 authors in the field of computer science, in this article we studied the relation between the centrality of authors in the coauthorship network and the future success of their publications. Clearly, there are certain limitations to our approach, which we discuss in the following.

First of all, any data-driven study of social behavior in general and citation behavior in particular is limited by the completeness and correctness of the used data set. The fact that name

³We use the R package *randomForest*, available at <http://cran.r-project.org/web/packages/randomForest/>

⁴For details on the procedure and the error estimates we refer to the Supplementary Material.

ambiguities are automatically resolved by the Microsoft Academic Search (MSAS) database by sophisticated and validated disambiguation heuristics is a clear advantage over simpler heuristics that have been used in similar studies.

In order to rule out effects that are due to different citation patterns in different disciplines, we limited our study to computer science, for which we expect the coverage of MSAS to be particular good. While this limits the generalization of our results to other fields, our work nevertheless represents – to the best of our knowledge – the first large-scale case study of social factors in citation practices. As publication practises seem to vary widely across disciplines, it will be interesting to investigate whether our results hold for other research communities as well.

Clearly, any study that tries to evaluate the *importance* or *centrality* of actors in a social network needs to be concerned about the choice of suitable centrality measures. In order to not overemphasize one particular – out of the many – dimensions of centrality in networks, we chose to use *complementary centrality measures* that capture different aspects of importance at the same time. The results of our prediction highlight that the combination of different measures is crucial – making clear that visibility and social influence are more complicated to capture than by a single centrality measures.

Finally, one may argue that our observation that authors with high centrality are cited more often is not a statement of a *direct causal relation* between centrality and citation numbers. After all, both centrality and citations could be secondary effects of, for instance, the scientific excellence of a particular researcher, which then translates into becoming central and highly cited at the same time. Clearly, we neither can – nor do we want – to rule out such possible explanations for our statistical findings. However, considering our finding of strong statistical dependence between social centrality and citation success, one could provocatively state the following: if citation-based measures were to be good proxies for scientific success, so should then be measures of centrality in the social network. We assume that not many researchers would approve having their work evaluated by means of such measures. We hence think that our findings are an important contribution to the ongoing debate about the meaningfulness and use of citation-based measures, as well as a better understanding of citation dynamics in general.

In summary, the contributions of our work are threefold:

1. We provide the, to the best of our knowledge, first large-scale study that analyses relations between the position of researchers in scientific collaboration networks and citation dynamics, using a set of complementary network-based centrality measures. A specific feature of our method is that we study *time-evolving* collaboration networks and citation numbers, thus allowing us to investigate possible mechanisms of social influence at a microscopic scale.

2. We show that – at least for the measures of centrality investigated in this paper – there is no *single* notion of centrality in social networks that could accurately predict the future citation success of an author. We expect this finding to be of interest for any general attempt to predict the success of actors based on their centrality in social networks.
3. Using modern machine learning techniques, we present a supervised classification method based on a Random Forest classifier, using a multidimensional feature vector of collaboration network centrality metrics. We show that this method allows for a remarkably precise prediction of the future citation success of a paper, solely based on the social embedding of its authors. With this, our method provides a clear indication for a strong statistical dependence between author centrality and citation success.

In conclusion, we provided evidence for a strong relation between the position of authors in scientific collaboration networks and their future success in terms of citations. We would like to emphasize that by this we *do not* want to join in the line of – sometimes remarkably uncritical – proponents of citation-based evaluation techniques. Instead, we hope to contribute to the discussion about the manifold influencing factors of citation measures and their explanatory power concerning scientific success. Especially, we *do not* see our contribution in the development of automated success prediction techniques, whose widespread adoption could possibly have devastating effects on the general scientific culture and attitude. Highlighting social influence mechanisms, we rather hope that our work contributes to a better understanding of the multifaceted, complex nature of citations, which should be a prerequisite for any reasonable application of citation-based measures.

7 Acknowledgement

EM, IS and FS acknowledge funding by the Swiss National Science Foundation, grant no. CR31I1_140644/1. AG acknowledges funding by the EU FET project MULTIPLEX 317532. We especially thank Microsoft Research for granting unrestricted access to the Microsoft Academic Search service.

References

- [1] D. E. Acuna, S. Allesina, and K. P. Kording. Future impact: Predicting scientific success. *Nature*, 489:201–202, September 2012.

- [2] M. Bras-Amorós, J. Domingo-Ferrer, and V. Torra. A bibliometric index based on the collaboration distance between cited and citing authors. *J. Informetrics*, 5(2):248–264, 2011.
- [3] L. Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
- [4] G. Csardi and T. Nepusz. The igraph software package for complex network research. *InterJournal, Complex Systems*, 1695(5), 2006.
- [5] W. D. Figg, L. Dunn, D. J. Liewehr, S. M. Steinberg, P. W. Thurman, J. C. Barrett, and J. Birkinshaw. Scientific collaboration results in higher citation rates of published articles. *Pharmacotherapy: The Journal of Human Pharmacology and Drug Therapy*, 26(6):759–767, 2006.
- [6] E. Garfield. Citation indexes for science: A new dimension in documentation through association of ideas. *Science*, 122(3159):108–111, 1955.
- [7] C. Goutte and E. Gaussier. A probabilistic interpretation of precision, recall and f-score, with implication for evaluation. In *Advances in Information Retrieval*, pages 345–359. Springer, 2005.
- [8] W. Heisenberg. *Der Teil und das Ganze: Gespräche im Umkreis der Atomphysik*. Piper und Co. Verlag München, 1969.
- [9] J. E. Hirsch. An index to quantify an individual’s scientific research output. *Proceedings of the National Academy of Sciences of the United States of America*, 102(46):16569–16572, 2005.
- [10] J. E. Hirsch. Does the h index have predictive power? *Proceedings of the National Academy of Sciences*, 104(49):19193–19198, 2007.
- [11] J. Hsu and D. Huang. Correlation between impact and collaboration. *Scientometrics*, 86(2):317–324, 2011.
- [12] J. Katz and D. Hicks. How much is a collaboration worth? a calibrated bibliometric model. *Scientometrics*, 40(3):541–554, 1997.
- [13] F. Laloë and R. Mosseri. Bibliometric evaluation of individual researchers: not even right... not even wrong! *Europhysics News*, 40(5):26–29, 2009.
- [14] L. Leydesdorff. Theories of citation? *Scientometrics*, 43(1):5–25, 1998.
- [15] H. B. Mann and D. B. Whitney. On a test of whether one of two random variables is stochastically larger than the other. *Annals of Mathematical Statistics*, 18(1):50–60, 1947.

- [16] T. Martin, B. Ball, B. Karrer, and M. E. J. Newman. Coauthorship and citation in scientific publishing. *arXiv preprint arXiv:1304.0473*, 2013.
- [17] M. E. J. Newman. The first-mover advantage in scientific publication. *EPL (Europhysics Letters)*, 86(6):68001, 2009.
- [18] M. E. J. Newman. *Networks: an introduction*. Oxford University Press, 2009.
- [19] M. E. J. Newman. Prediction of highly cited papers. *ArXiv preprint arXiv:1310.8220*, Oct. 2013.
- [20] J. Nicolaisen. The social act of citing: Towards new horizons in citation theory. *Proceedings of the American Society for Information Science and Technology*, 40(1):12–20, 2003.
- [21] F. Radicchi, S. Fortunato, and C. Castellano. Universality of citation distributions: Toward an objective measure of scientific impact. *Proceedings of the National Academy of Sciences*, 105(45):17268–17272, 2008.
- [22] M. J. Stringer, M. Sales-Pardo, and L. A. N. Amaral. Statistical validation of a global model for the distribution of the ultimate number of citations accrued by papers published in a scientific journal. *Journal of the American Society for Information Science and Technology*, 61(7):1377–1385, 2010.
- [23] M. L. Wallace, V. Lariviere, and Y. Gingras. A small world of citations? the influence of collaboration networks on citation practices. *PLoS ONE*, 7(3):e33339, 03 2012.
- [24] S. Wasserman and K. Faust. *Social Network Analysis: Methods and Applications*. Cambridge University Press, 1994.

Supplementary Material

Centrality Metrics

There are many network metrics that can be used for social network analysis [18, 24]. Here we are quickly going to review the metrics we have been using in this work and their interpretation in coauthorship networks.

Degree Centrality The degree centrality of a node is its number of first-order neighbors, i.e. the number of nodes this node connects to via one link. In a directed network, this measure is divided into an in-degree centrality and an out-degree centrality. Since the here considered coauthorship network is undirected, the degree centrality is simply the number of its direct neighbors. Degree centrality is a local measure, as it does not depend on any global network properties other than the number of its neighbors. In the coauthorship network the degree centrality of a node is its number of coauthors.

Eigenvector Centrality The eigenvector centrality of a node is a global centrality measure, as non-local changes in the network can alter the node's eigenvector centrality. In short, a node has high eigenvector centrality if it is connected to other nodes with high eigenvector centrality. As such, this centrality measure goes beyond degree centrality as a mere measure of quantity (the number of neighbors) in that it introduces a notion of inheritance of importance. Used often, especially in its variant PageRank, eigenvector centrality of node v is the v th component of the Perron-Frobenius-eigenvector of the network's adjacency matrix. In the coauthorship network, eigenvector centrality has a meaning of importance, if one assumes that an author is more important if she coauthors papers with other authors of high importance.

Betweenness Centrality The Betweenness centrality is another often used global centrality measure. A node has high betweenness centrality if it lies on many shortest paths of the network. Hence, this centrality measure is a measure of importance in terms of network flows. If a node with high betweenness centrality would be removed, a lot of network flows would become less efficient, as the average length of shortest paths will increase. In the coauthorship network, a node with high betweenness centrality could be interpreted as a node with high importance for "fast knowledge transfer", as this person lies on many shortest paths connecting authors and their research.

K-Core Centrality The k -core centrality is a global centrality measure thought to measure the "coreness" of a node, i.e. how deep it is embedded in the network. A node has k -core centrality k if, when consecutively removing nodes that have degree 1, 2, ... $k - 1$ from the network, this node has not been removed, but will be removed in a next step when nodes with degree k are removed. k -core centrality is somewhat similar to eigenvector centrality, as a node must have neighbors with high k -core in order to have high-core itself. Different from eigenvector centrality,

k -core centrality is not additive. Hence compensating a low number of high k -core neighbors with a high number of low k -core neighbors does not guarantee the node to have high k -core. In the coauthorship network a node has a high k -core, if it is connected to many nodes that have high k -core themselves.

Correlations Between Citation Numbers and Centrality Metrics

In section *Effects of Author Centrality on Citation Success* of the main manuscript, we argue that citation numbers of an article, five years after its publication, are not Pearson- and Spearman-correlated with social network centrality metrics of its authors. Table 6 summarizes the Pearson and Spearman correlation coefficients for the considered metrics. None of these results allows to conclude any significant correlation.

	Pearson r	Spearman ρ
k-core	0.05	0.15
Eigenvector	0.01	0.07
Betweenness	0.13	0.15
Degree	0.05	0.16

Table 6: Pearson and Spearman coefficients measuring correlations between citation numbers of a paper (five years after publication) and coauthorship network centrality of its authors.

Precision and Recall

In Machine Learning it is standard practice to assess the goodness of a classifier using the quantities *precision* and *recall* [7].

Precision is defined as

$$Precision = \frac{TruePositives}{TruePositives + FalsePositives}. \quad (1)$$

It hence is equal to the fraction of correctly predicted instances compared to all predicted instances. As such, precision quantifies how reliable the predicted results are. However, it does not make any statement about how many relevant results the predictor returns. For example, a simple predictor could be to always return one element from which it is known ex ante, that it is a true prediction. In this scenario precision would be 100%, however the sensitivity might be poor as there might be more than one relevant element. This last point is quantified using *recall*.

Recall (or *sensitivity*) is defined as

$$Recall = \frac{TruePositives}{Positives}. \quad (2)$$

In a case of 100 samples, 1 positive and 99 negative, a classifier that solely returns the one positive element has *recall* = 1 as well as *precision* = 1. However, a classifier that returns all 100 elements still has *recall* = 1, although precision will be *precision* = 1/99.

Precision and recall are usually presented jointly to assess the goodness of a classifier. Sometimes they are combined to the so called *F-score*:

$$F = 2 \cdot \frac{precision \cdot recall}{precision + recall} \quad (3)$$

F-score value $F = 1$ indicates best possible classifier goodness, $F = 0$ indicates worst.

Random Forest Prediction

In section four *Predicting Successful Publications* of the main manuscript, we use a Random Forests classifier [3] to predict success of publications based on coauthorship centrality of its authors. A Random Forest classifier fits a number of classification trees (a so called forest) on bootstrap samples of the data set with subsequent averaging over all individual classification trees to improve predictive accuracy. Random Forests perform well on classification tasks involving a large number of features. We will quickly outline the high-level procedure, for details please see [3].

1. A bootstrap sample from the minority class, and a sample with the same number of cases from the majority class are drawn, with replacement, for each iteration of the Random Forest classification.
2. For each iteration, a classification tree from the data is grown, without pruning. Assuming that there are M variables in the data set, a number $m \ll M$ is chosen such that at each node, m variables are selected at random out of the M and the best split on these m is used to split the node. The best split is discovered using the CART algorithm, where, at each node, instead of searching through all variables for the optimal split, only m variables are searched through.
3. The two steps above are repeated depending on how many trees are desired to be grown. After a bootstrap sample of the data is put down each tree, each tree votes for what it computes as the true classes for each case in that sample set. The final classification for each case is the one having most votes among the individual trees.

In this work, we utilize the R-package *randomForest*⁵ to perform the classification.

Out-of-bag (OOB) Error Estimate

Random Forests do not require a separate test set to cross-validate the classification and to estimate the classification error rate. Instead, error estimation can be done during run time of the algorithm. Each tree is constructed using a different bootstrap sample from the original data where about one-third of the cases from the sample are not used in the construction. These are called the *Out of bag (OOB)* cases. These *OOB* cases are later classified using the constructed decision trees. At the end of the classification, an error estimate is computed by averaging the proportion of times where the class that got the most votes is not equal to the true class over all *OOB* cases.

Importance and Significance of Features

In addition to OOB error estimation, Random Forest classifiers quantify the significance and importance of features.

1. For every tree in the forest which uses a sample with m features, the number of votes for the correct classification of the *OOB* cases are calculated.
2. Then the values of feature k in the *OOB* cases are permuted, and the correctly classified cases from this permutation is also calculated.
3. The average of the difference of counts from step 1 and step 2 over all trees is the importance score of feature k .
4. If the values of this score are independent across individual trees, the standard error can be computed by dividing the raw score by its standard error to get a *z-score* and assigning a significance level to the feature.

Taking samples of three variables at a time ($m = 3$) for each split, we classify the publications. Using the internal variable importance ranks, we are able to detect which variables hold the most predictive power for classifying true positives.

⁵<http://cran.r-project.org/web/packages/randomForest/>